

NAME	DATE

Bias- Variance tradeoff

Train set error	1%	15%	15%	0.5%
Dev set error	11%	16%	30%	1%

high high high low
variance bias bias bias

if train set error is large → high bias

if diff. between train set error & dev set error is large → high variance

Above all is true, if optimal error / Bayes error (Human error) is ≈ 0% & train set & dev set is drawn from the same distribution

Basic recipe for ML

Try these stuffs

High bias \rightarrow Use bigger network
 \rightarrow Train larger data
 (Try new NN architecture)

High variance \rightarrow More data

\rightarrow Regulariz?

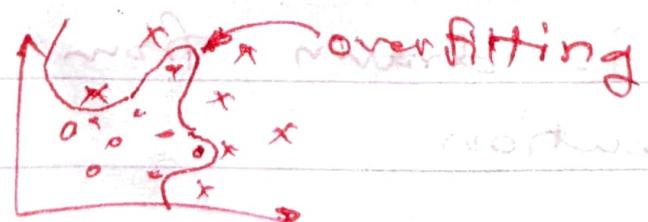
(Try new NN architecture)

bias \rightarrow error your model has

It fails while capturing some

trend & the features in data

Variance \rightarrow overfits the data



Regularization

Frobenius norm formula

$$\|w^{[l]}\|^2 = \sum_{i=1}^n \sum_{j=1}^{n^{[l-1]}} (w_{i,j}^{[l]})^2$$

i of the matrix \rightarrow no. of neurons
of its layer

j of the matrix \rightarrow no. of neurons
in the previous layer $n^{[l-1]}$.

[Hence, summation from $j=1$ to $n^{[l-1]}$]

$$dw^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} \cdot w^{[l]}$$

(new formula) for gradient

$$w^{[l]} := w^{[l]} - \alpha dw$$

$$w^{[l]} = w^{[l]} - \alpha [(\text{from backprop}) \\ \text{prev value} + \frac{\lambda}{m} \cdot w^{[l]}]$$

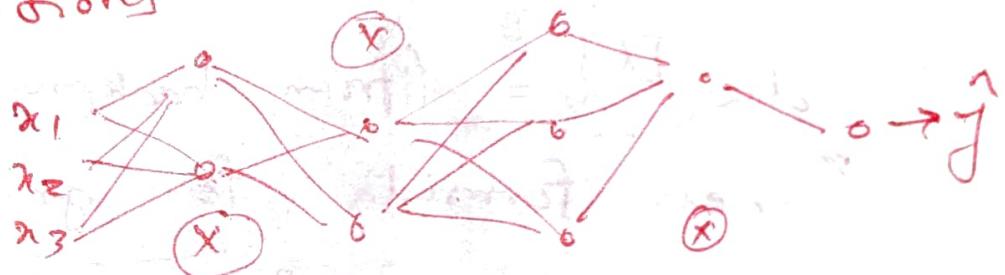
$$(1 - \frac{\alpha \lambda}{m}) w^{[l]} = w^{[l]} - \frac{\alpha \lambda}{m} [l] \\ - \alpha (\text{from backprop})$$

Dropout Regularization

Inverted dropout (Most popular version of dropout)

Turning some of the neurons in NN to zero (dropping them off)
OR deactivating them.

We can do this by simply setting value of $w_{i,j}^{[l]}$ to zero rather than changing $w^{[l]}$'s dimensions



Say we want to reduce number

(of neurons in layer 3 by 20%.

To keep 80% of the neurons alive

`layer3 = np.random.rand(a3.shape[0], a3.shape[1]) < keep_prob`

$$\text{keep_prob} = 0.8$$

~~a^3~~ = np. multiply (a^3 , d^3)

$a^3 = \text{np. multiply } (a^3, d^3) \text{ & multiply}$
 OR $a^3 * d^3 = d^3$ Jing by
 d^3

✓ d^3 is a matrix of 0's & 1's

d^3 is a matrix of 0's & 1's

$a^3 / \text{keep_prob}$

Why divide a^3 by keep_prob?

When you will multiply a^3 with d^3 , some of the neurons will be set to zero. To compensate this, we are stepping up the remaining neurons by multiplying it with keep_prob. This done to preserve the expected value of

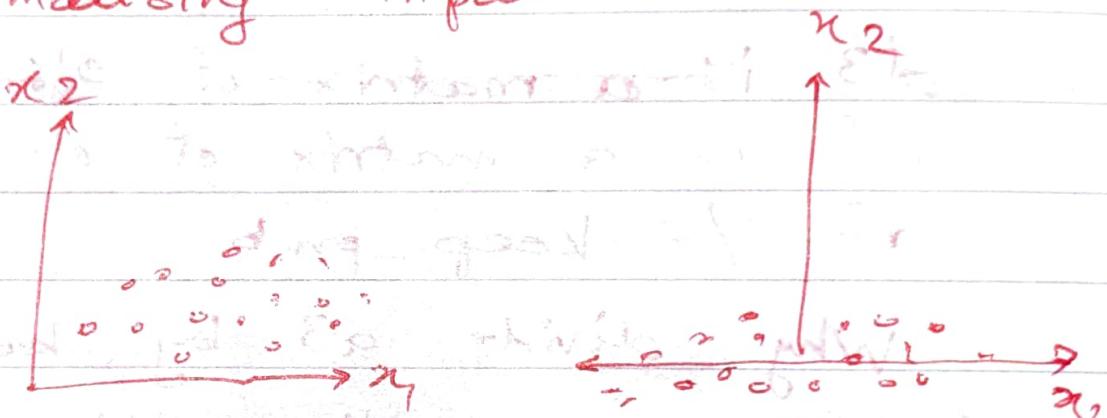
$z^{[4]}$

$$z^{[4]} = w \cdot a^{[3]} + b^{[4]}$$

Hence, this is called inverted dropout

do this only for training set not test set. You can tune value of keep-prob

Normalizing Input



subtract mean

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x' = x - \mu$$

x_2'

(3)

x_1' has high variance than x_2' . We can

normalize this

variance by performing

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i) \times 2}$$

classifying parallel lines / parallel
or to no classification of the lines)

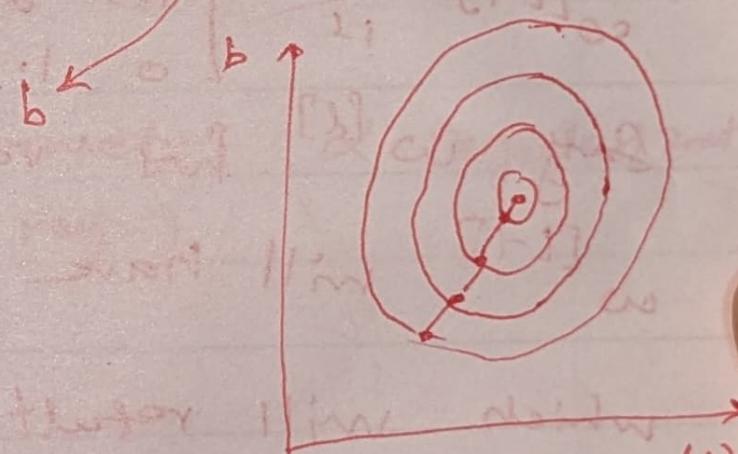
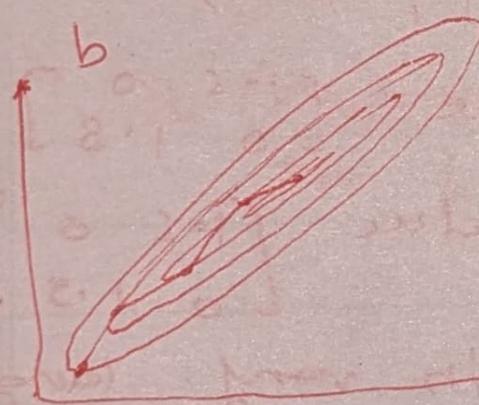
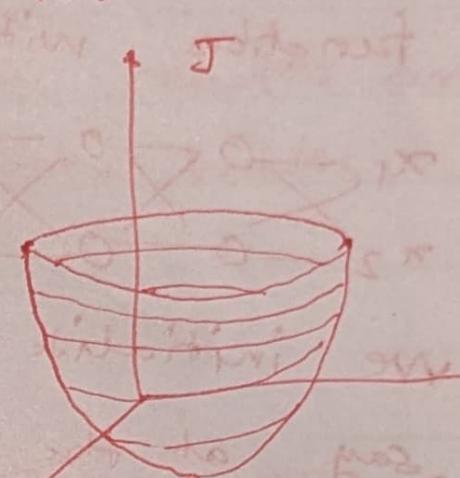
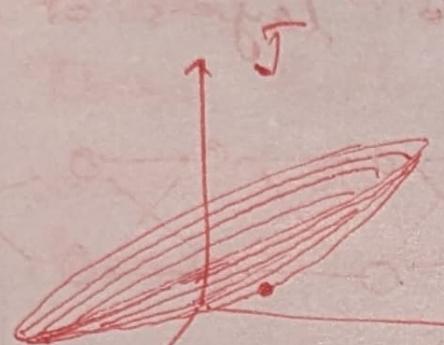
info from one input works well & if

$\alpha = \alpha/6$ working with (1)

so Why normalize inputs?

Unnormalised

Normalized

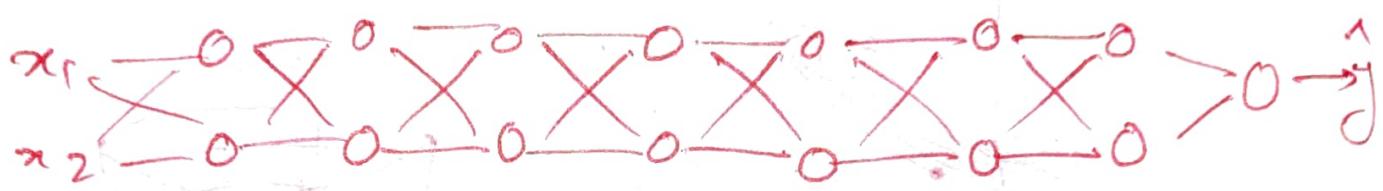


taking large value of w kind of any
& might affect learning value of w
as gradient descent as you can
will take a lot of steps start from any-
to oscillate back & forth to find min. when.

Vanishing / exploding gradients
Caused due to initialization of w .

In a NN where there are many layers (L), this problem may arise.

Say we have 2 features x_1 & x_2 and all the hidden also has 2 activation functions with 8-10 layers of it.



We initialize b with 0.
~~say at the output layer, value of $c^{[L-1]}$ is $\begin{bmatrix} 1.5 & 0 \end{bmatrix}$, then~~

Say $w^{[L]}$ has value $\begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$ so
 $w^{[L-1]}$ will have value $\begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1}$

which will result in very large values of $w^{[L-1]}$ and y .

Gradient will explode.

Or if ω has value $\begin{bmatrix} 0.9 & 0 \\ 0 & 0.9 \end{bmatrix}$

or something less 1, gradient descent will vanish.

These Both of those results are undesirable bcoz these make learning slow

keep - prob α regularization

λ (hyperparameter α regularization
in L2 reg.)

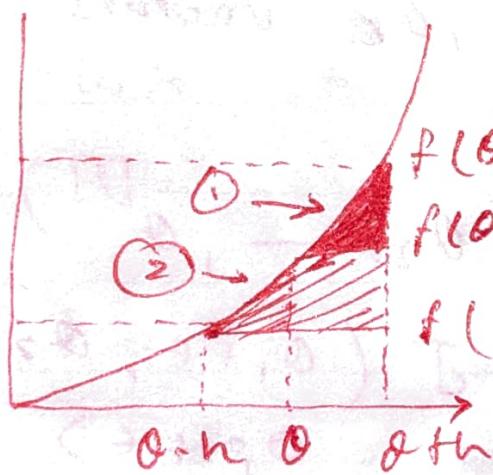
Numerical approximation of Gradients

We often know

$$\frac{df(x)}{dx} = \lim_{n \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$
$$= \lim_{x \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

But, using $\lim_{n \rightarrow 0} \frac{f(x+h) - f(x-h)}{2h}$

This is more accurate.



say $f(\theta)$ is the
function $= \theta^3$
previously
we used only
triangle ①

But we can use ① & ② to
get more accurate results

$$\Delta f'(\theta) = \lim_{h \rightarrow 0} \frac{f(\theta+h) - f(\theta-h)}{2h}$$

If we want to check the computed gradients i.e. d_w , db , d_a , d_b , we can use the approximation technique.

We know J is function of

$$(w^1, b^1), (w^2, b^2), \dots, (w^{[L]}, b^{[L]})$$

OR $\theta_1, \theta_2, \dots, \theta^{[L]}$.

what we can do is merge all the $\theta_1, \theta_2, \dots, \theta^{[L]}$ vectors into single vector θ

$$J(\theta) = J(\theta_1, \theta_2, \dots, \theta^{[L]})$$

$$d\theta_{\text{approx.}}[\ell] = J(\theta_1 + \epsilon, \theta_2 + \epsilon, \dots, \theta^{[\ell]} + \epsilon) -$$

$$\frac{J(\theta_1, \theta_2, \dots, \theta^{[\ell]} - \epsilon) - J(\theta_1, \theta_2, \dots, \theta^{[\ell]} + \epsilon)}{2\epsilon}$$

$$\frac{(J(\theta_1, \theta_2, \dots, \theta^{[\ell]} + \epsilon) - J(\theta_1, \theta_2, \dots, \theta^{[\ell]} - \epsilon))}{2\epsilon}$$

$$\approx d\theta^{[\ell]} = \frac{\partial J}{\partial \theta^{[\ell]}}$$

$$\|d\theta_{approx} \cdot \text{shape}\| = \|d\theta_{approx}\| \cdot \|\text{shape}\|$$

check

$$\text{on L2 norm} \rightarrow \frac{\|d\theta_{approx} + d\theta\|_2}{\|d\theta\|_2}$$

$$\|d\theta_{approx}\|_2 + \|d\theta\|_2$$

Usually, $\varepsilon = 10^{-3}$

if the value of this is

close to 10^{-3} \rightarrow good res.

10^{-5} \rightarrow bad res.

10^{-3} \rightarrow wrong

Weigh Initialization

Instead of randomly generating weights, we can initialize them to a particular value for better performance.

$$W^{[l]} = np.random.rand(\text{shape})$$

$$\text{where } n^{[l]} = \sqrt{n^{[l-1]}} \left(\frac{1}{\epsilon} \right)$$

$n^{[l]}$ number of neurons in last layer

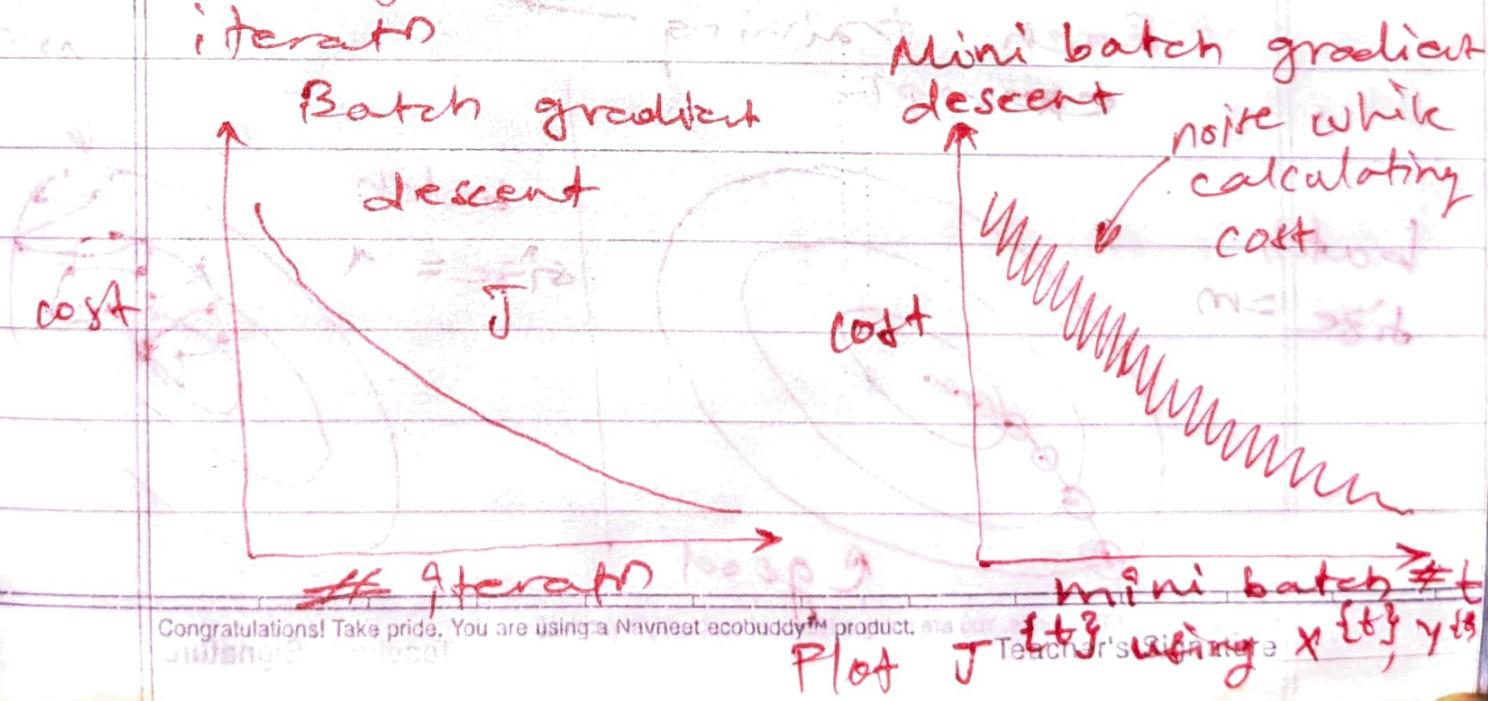
$$\text{For ReLU, it is } np.sqrt\left(\frac{2}{n^{[l-1]}}\right)$$

$$\text{For tanh, it is } np.sqrt\left(\frac{1}{n^{[l-1]}}\right)$$

Mini batch gradient descent

Instead ~~of~~ of training on entire training set, we'll divide it into 't' parts. Then train ^{model on} each t part.

say we divide into 1000 examples each. We'll perform forward prop on first 1000 examples, then perform back prop on those 1000 examples by calculating cost & g.d.'s for number of iterations.



Mini-batch gd is faster & computationally cheaper

• Mini-batch size m

• if batch size = 1, then it's SGD

if batch size = m (then good

for OR same as batch gd.)

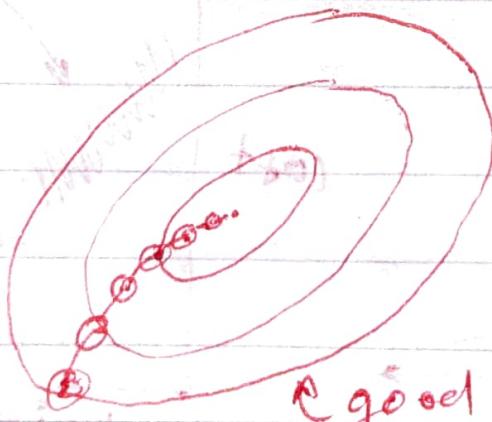
• if batch size = 1 & size

{ (Stochastic gd)

$$(x^{(t)}, y^{(t)}) = [(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})]$$

Each training example

batch
size = m



good

batch
size = 1



Using batch size of n is same as using for loops, for training the model & not taking advantage of vectorization.

A mini batch size should be b < m . [m = No. of training eg]

If $m \leq 2000$, don't do mini batch

Typical mini batch size:

64, 128, 256, 512

2^6 2^7 2^8 2^9

power of 2^k

Now mini-batch is also a hyperparameter

Exponentially weighted avg.

$$V_t = \beta V_{t-1} + (1-\beta) \bar{o}_t$$

avg of t^{th} time / example

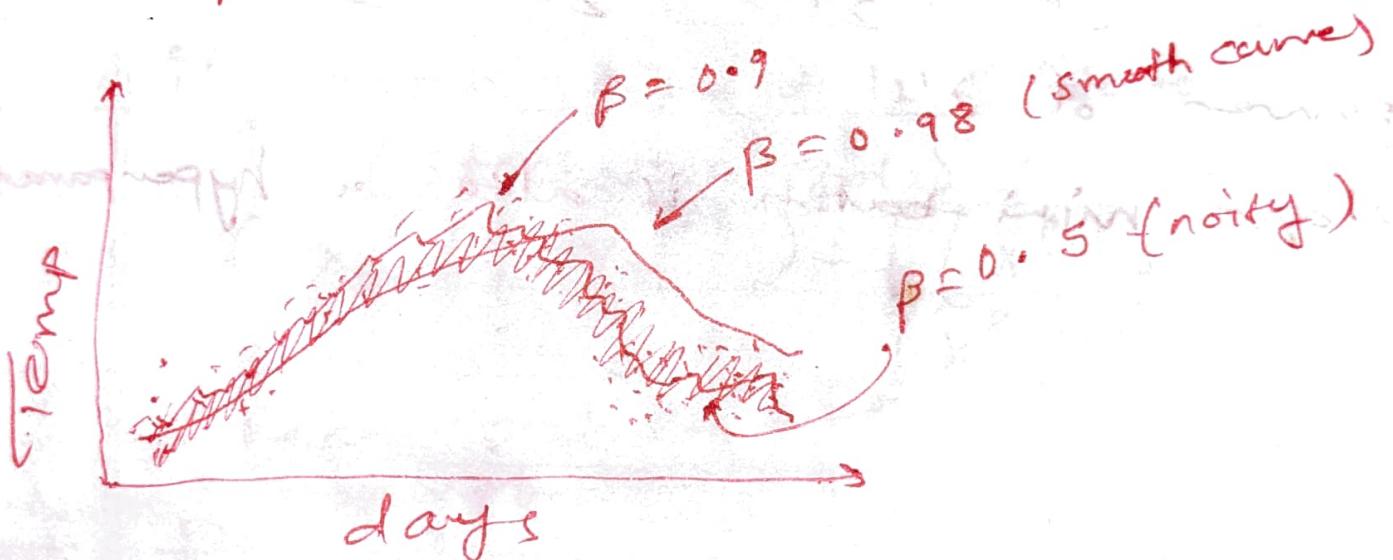
[V_t as approximately averaging

over $\frac{1}{1-\beta}$ days

if $\beta = 0.9$, V_t is avg. over 10 days

$\beta = 0.98$, 50 days

$\beta = 0.5$, 2 days



β is also a hyperparameter

Gradient descent with momentum
add's a term to speed up process of gd.

On iteration t:

- Compute dW , db on the current mini-batch
- momentum \rightarrow acceleration

$$Vdw = \beta Vdw + (1-\beta) dW$$

$$Vdb = \beta Vdb + (1-\beta) db$$

$$w = w - \alpha Vdw, b = b - \alpha Vdb$$



Analogy: ball rolling down hill

Hyperparameters: α , β

RMS prop [Root mean squared prop]

On iteratⁿ t:

Compute dW , db on curr. mini-batch

$$Sdw = \beta Sdw + (1-\beta) dW^2 \leftarrow \text{small element wise}$$

$$Sdb = \beta Sdb + (1-\beta) db^2 \leftarrow \text{large}$$

$$\hat{w} = \frac{w - \alpha dW}{\sqrt{Sdw}}, \hat{b} = \frac{b - \alpha db}{\sqrt{Sdb}}$$

Adam optimization & algorithm

$$V_{dw} = 0, S_{dw} = 0 \quad V_{db} = 0, S_{db} = 0$$

In Adam, we'll be using doing bias correction

$$V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1^t), V_{db}^{\text{corr.}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{\text{corr.}} = S_{dw} / (1 - \beta_2^t), S_{db}^{\text{corr.}} = S_{db} / (1 - \beta_2^t)$$

$$\hat{w} := w - \alpha \cdot \frac{V_{dw}^{\text{corr.}}}{\sqrt{S_{dw}^{\text{corr.}}} + \epsilon}, \hat{b} := b - \alpha \cdot \frac{V_{db}^{\text{corr.}}}{\sqrt{S_{db}^{\text{corr.}}} + \epsilon}$$

Why $\sqrt{S + \epsilon}$ in denom?

→ if $S \approx 0$, the \sqrt{S} term will be very large to prevent this ϵ is added

$$\epsilon \approx 10^{-8}$$

α : needs to tuned

$$\beta_1 := 0.9$$

$$\beta_2 := 0.999$$

$$\epsilon := 10^{-8}$$

} Hyper-parameters

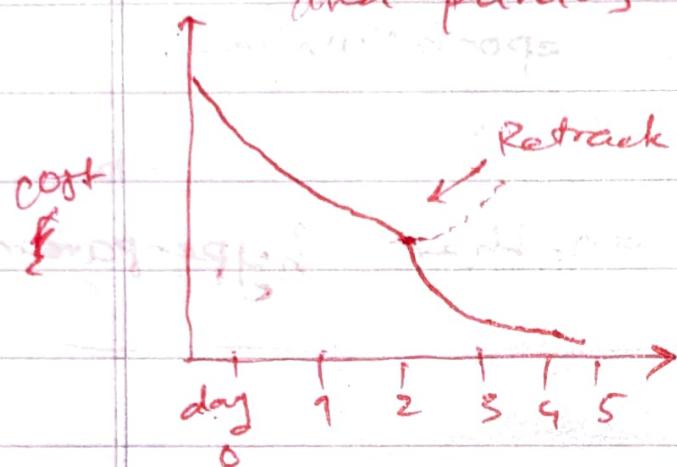
Learning rate decay (α)

$$\frac{\alpha_0}{1 + \text{decayRate} \times \text{epoch Number}}$$

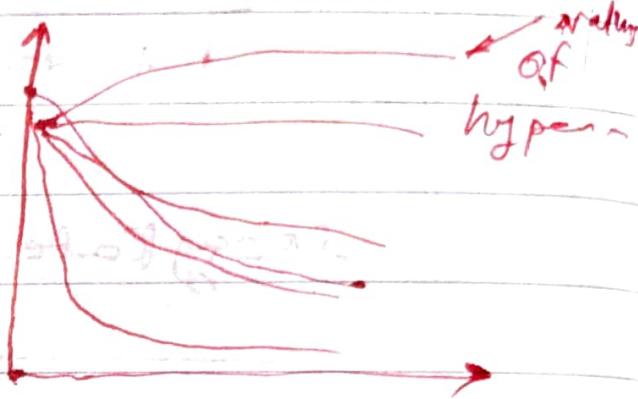
decayRate is another hyperparameter

Training models by ~~based~~ changing hyperparameters

aka pandas



aka Caviar diff. values of hyper-



Training the model

→ by changing the values of hyperparameters day-by-day to get best results

Time consuming

but save lot of computational resources

Train the many

models in parallel with each having diff. values of hyperparameters

fast but takes heavy toll of computational resources

Normalizing activations in a network
This makes training process smoother & efficient

Previously we used to perform normalization on input layer / features. Now, if we perform it on all layers, we would get better results. This is known as

'Batch normalization'

$$\text{Mean } \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x = x - \mu$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m [x^{(i)} - \mu]^2 \leftarrow \text{element-wise}$$

For $x = x / \sigma^2$

input layer

For l layers: $\underbrace{z^{(1)}, \dots, z^{(l)}}_{z^{[l]}}$, $\overbrace{z^{(l+1)}, \dots, z^{(m)}}^{z^{[m-l]}}$

for some
 $\ell^{\text{th}} \rightarrow \mu = \frac{1}{m} \sum_i z^{(\ell)}$

$$\sigma^2 = \frac{1}{m} \sum_i (z^{(\ell)} - \mu)^2$$

$$z_{\text{norm}}^{(\ell)} = \frac{z^{(\ell)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$z^{(\ell)} = \gamma z_{\text{norm}}^{(\ell)} + \beta$$

Use \hat{z} instead of $z^{(\ell)}$

γ & β are learnable parameters

Batch norm as regularization

- Batch norm is applied on each mini batch. Mini-batch has size typically ranging from 64 to 512 (very small compared to all the training examples).
- Each mini batch is scaled by the mean / variance computed on just that mini-batch.
- This adds noise to the values $z^{[l]}$ within that minibatch. Similar to dropout, it adds noise to each hidden layer's activation.
- This has a slight regularization effect.

Dropout property:

Size of minibatch: $64 \rightarrow 512$
more \rightarrow less regularization
 $(1) 64 \rightarrow 512$ regularization

Softmax (for C classes as output)

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{i=1}^C t_i}$$

say $z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$

$$a^{[L]} = g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5 / (e^5 + e^2 + e^{-1} + e^3) \\ e^2 / (e^5 + e^2 + e^{-1} + e^3) \\ e^{-1} / (e^5 + e^2 + e^{-1} + e^3) \\ e^3 / (e^5 + e^2 + e^{-1} + e^3) \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \leftarrow \text{hard max}$$

$$\begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \leftarrow \text{softmax}$$

Then loss function becomes :

$$L(\hat{y}, y) = - \sum_{j=1}^4 y_j \log \hat{y}_j$$

$$\text{cost} = J(w^{[1]}, b^{[1]}, \dots) = \overline{\sum_{i=1}^m L(\hat{y}_i, y_i)}$$

$$Y = [y^{[1]}, y^{[2]}, \dots, y^{[m]}]$$

$$= \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix} \quad \text{Dimensions : } 4 \times m$$